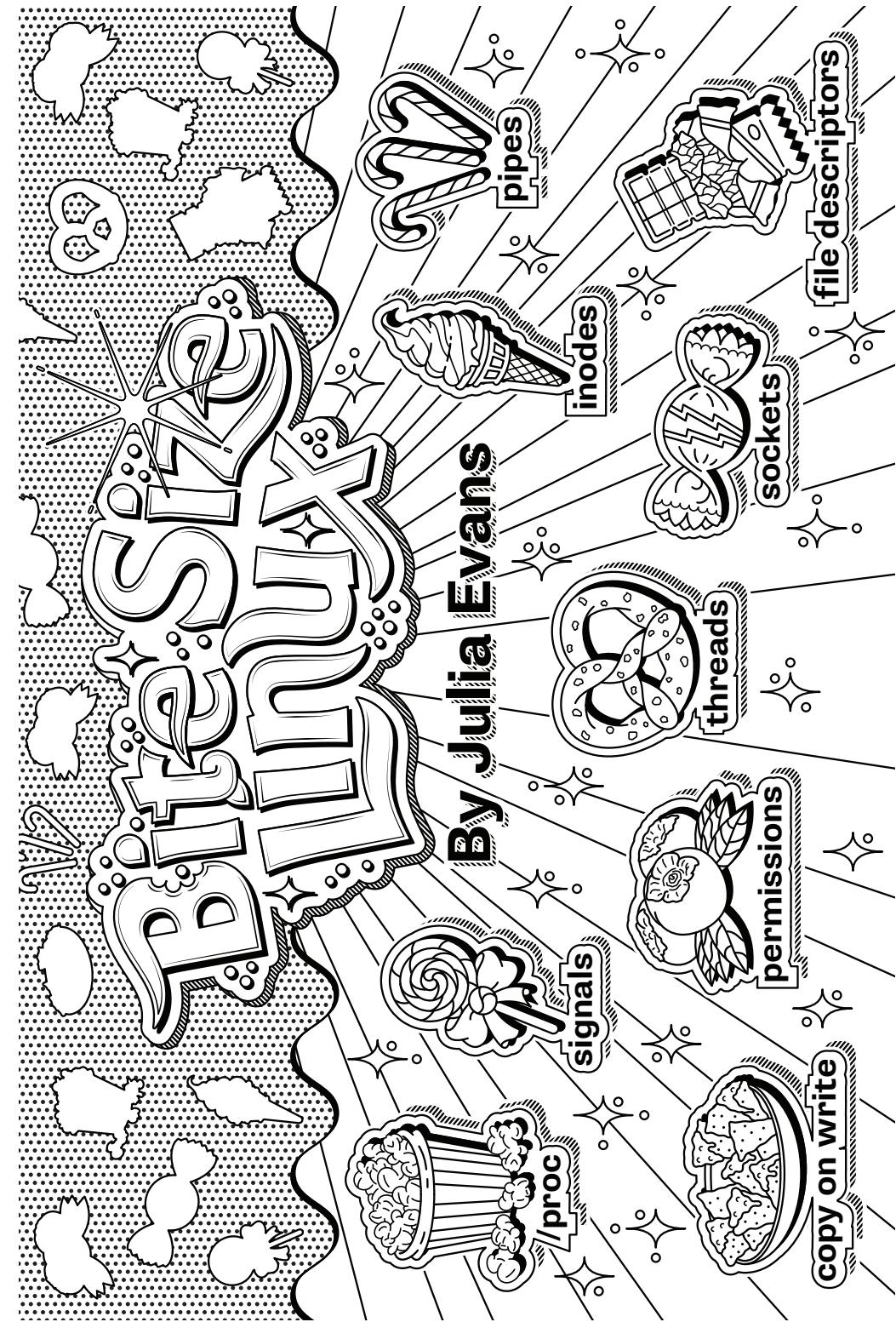


I love this?

find more awesome zines at  
→ [wizardzines.com](http://wizardzines.com) ←





I want to learn more?  
I highly recommend  
this book: →

Every chapter is a readable,  
short (usually 10-20 pages)  
explanation of a Linux system.  
I used it as a reference  
constantly when writing  
this zine.

I & it because even though  
it's huge and comprehensive  
(1500 pages!), the chapters  
are short and self-contained  
and it's very easy to pick it  
up and learn something.

by Julia Evans  
<https://jvns.ca>  
[twitter.com/b0rk](https://twitter.com/b0rk)

... 5 minutes later ...

oh wow! That's  
not so complicated.  
I want to learn  
more now!

you're in the  
right place! This  
zine has 19 comics  
explaining important  
Linux concepts.

Julia

# man page sections 22

man pages are split up into 8 sections	man page sections
① ② ③ ④ ⑤ ⑥ ⑦ ⑧	① programs \$ man grep \$ man ls
\$ man 2 read	② system calls \$ man sendfile \$ man ptrace
means "get me the man page for read from section 2".	③ C functions \$ man printf \$ man fopen
There's both → a program called "read" → and a system call called "read"	④ devices \$ man null for /dev/null docs
so \$ man 1 read	⑤ file formats \$ man sudoers for /etc/sudoers
gives you a different man page from	⑥ games \$ man sl \$ man proc ! files in /proc !
\$ man 2 read	⑦ miscellaneous explains concepts! \$ man '7 pipe \$ man 7 symlink
If you don't specify a section, man will look through all the sections & show the first one it finds.	⑧ sysadmin programs \$ man apt \$ man chroot

## ◆ Table of contents ◆

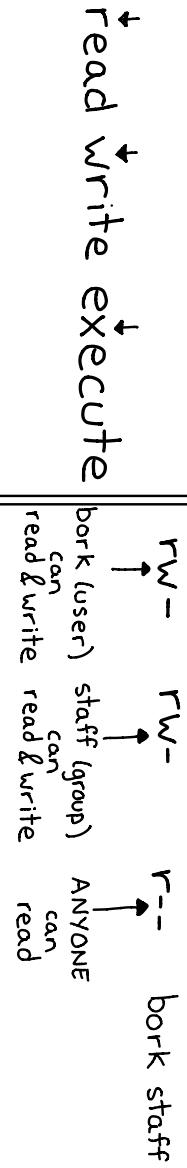
unix permissions...4	sockets.....10	virtual memory...17
/proc .....5	unix domain sockets.....11	shared libraries...18
system calls.....6	processes.....12	copy on write ....19
signals.....7	threads.....13	page faults.....20
file descriptors.8	floating point....14	mmap.....21
pipes.....9	file buffering .....15	man page sections.....22
		memory allocation.....16

# Unix permissions

4

There are 3 things you can do to a file

↓ read ↓ write execute



File permissions are 12 bits  
 setuid setgid  
 000 user group all  
 sticky rwx rwx rwx  
 For files:  
 r = can read  
 w = can write  
 x = can execute  
 For directories, it's approximately:  
 r = can list files  
 w = can create files  
 x = can cd into & access files

110 in binary is 6  
 So rw- r-- r--  
 = 110 100 100  
 = 6 4 4  
 chmod 644 file.txt  
 means change the permissions to:  
 rwx - r-- r--  
 Simple!

setuid affects executables  
 \$ ls -l /bin/ping  
 rws r-x r-x root root  
 this means ping always runs as root  
 setgid does 3 different unrelated things for executables, directories, and regular files.  
 unix! story it's a long story

# mmap

21

What's mmap for?  
 I want to work with a VERY LARGE FILE but it won't fit in memory  
 You could try mmap!  
 (mmap = "memory map")

load files lazily with mmap  
 When you mmap a file, it gets mapped into your program's memory.

2TB file → 2TB of virtual memory  
 but nothing is ACTUALLY read into RAM until you try to access the memory.  
 (how it works: page faults!)

how to mmap in Python  
 import mmap  
 f = open("HUGE.txt")  
 mm = mmap.mmap(f.fileno(), 0)  
 ↗ this won't read the file from disk!  
 Finishes ~instantly.

print(mm[1000:1])  
 this ↑ will read only the last 1000 bytes!

sharing big files with mmap  
 we all want to read the same file!  
 no problem! mmap

dynamic linking uses mmap  
 I need to use libc.so.6  
 ↗ standard library

Even if 10 processes mmap a file, it will only be read into memory once!

you too eh? no problem.  
 I always mmap, so that file is probably dynamic linker already.

# page faults

20

every Linux process has a page table

## \* page table \*

virtual memory | physical memory address

0x19723000	0x1420000
0x19724000	0x1423000
0x1524000	not in memory
0x1844000	0x4a000 read only

some pages are marked as either

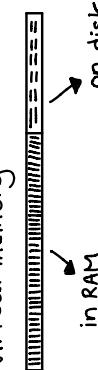
- ★ read only

★ not resident in memory

when you try to access a page that's marked "not resident in memory", it triggers a ! page fault!

"not resident in memory" usually means the data is on disk!

virtual memory



Having some virtual memory that is actually on disk is how swap and mmap work.

What happens during a page fault?

- the MMU sends an interrupt
- your program stops running
- Linux kernel code to handle the page fault runs

I'll fix the problem and let your program keep running

## how swap works

- ① run out of RAM  

A diagram showing a vertical stack of memory blocks labeled "RAM" overflowing into a horizontal bar labeled "disk".
- ② Linux saves some RAM data to disk  

A diagram showing a vertical stack of memory blocks labeled "RAM" with an arrow pointing down to a horizontal bar labeled "disk".
- ③ mark those pages as "not resident in memory" in the page table  

A diagram showing a vertical stack of memory blocks labeled "virtual memory" with an arrow pointing down to a horizontal bar labeled "RAM". The "virtual memory" stack has several blocks crossed out with a diagonal line, while the "RAM" stack remains intact.
- ④ when a program tries to access the memory, there's a ! page fault!  

A diagram showing a vertical stack of memory blocks labeled "virtual memory" with an arrow pointing down to a horizontal bar labeled "Linux". A speech bubble from the "Linux" bar says "time to move some data back to RAM!" with an arrow pointing to the "RAM" bar.
- ⑤ time to move some data back to RAM!  

A diagram showing a vertical stack of memory blocks labeled "virtual memory" with an arrow pointing down to a horizontal bar labeled "RAM". The "virtual memory" stack has several blocks crossed out with a diagonal line, while the "RAM" stack remains intact.

# an amazing directory: /proc 5

Every process on Linux has a PID (process ID) like 42.

In /proc/42, there's a lot of VERY USEFUL information about process 42.

/proc/PID/cmdline command line arguments the process was started with

/proc/PID/environ all of the process's environment variables

/proc/PID/stack

The kernel's current stack for the process. Useful if it's stuck in a system call.

/proc/PID/maps

Directory with every file the process has open! Run \$ ls -l /proc/42/fd to see the list of files for process 42. These symlinks are also magic & you can use them to recover deleted files!

/proc/PID/exe symlink to the process's binary magic: works even if the binary has been deleted!

/proc/PID/status Is the program running or asleep? How much memory is it using? And much more!

and `more`

Look at

man proc

for more information!

# System calls

6

The Linux kernel has code to do a lot of things

- read from a hard drive
- make network connections
- create new process
- kill a process
- change file permissions
- keyboards
- change file drivers

every program uses system calls

`(:)` I use the 'open' syscall to open files

`(:)` Python program  
`(:)` me too!

`(:)` Java program  
`(:)` me three!

`(:)` C program

and every system call has a number (e.g. chmod is #90 on x86-64)

So what's actually going on when you change a file's permissions is:

`(:)` run Syscall #90 with these arguments

`(:)` ok! Linux

you can see which system calls a program is using with `{strace}`

\$ strace ls /tmp

will show you every system call 'ls' uses!

it's really fun!

so don't run it on your production database

## COPY ON WRITE

19

On Linux, you start new processes using the `fork()` or `clone()` system call.

calling fork creates a child process that's a copy of the caller

`(:)` Parent  
`(:)` child

so Linux lets them share physical RAM and only copies the memory when one of them tries to write

`(:)` I'd like to change that memory

`(:)` okay! I'll make you your own copy!

`(:)` Linux

the cloned process has EXACTLY the same memory.

- same heap
- same stack
- same memory maps

if the parent has 3GB of memory, the child will too.

copying all that memory every time we fork would be slow and a waste of RAM

often processes call exec right after fork, which means they don't use the parent process's memory basically at all!

When a process tries to write to a shared memory address:

- ① there's a `page fault`
- ② Linux makes a copy of the page & updates the page table
- ③ the process continues, blissfully ignorant

`(:)` .. It's just like I have my own copy

# shared libraries

18

Most programs on Linux use a bunch of C libraries:

openSSL  
(for SSL)  
sqlite  
(embedded db!)  
lib pcre  
(regular  
expressions!)  
z lib  
(gzip!)  
libstdc++  
(C++ standard library!)

There are 2 ways to use any library:

- ① Link it into your binary  
`your code` `z lib` `sqlite`
- ② Use separate shared libraries  
`your code` `z lib` `sqlite` ↗ all different files

how can I tell what shared libraries a program is using?

`ldd` /usr/bin/curl  
libz.so.1 => /lib/x86\_64...  
libresolv.so.2 =>...  
libc.so.6 =>...  
+ 34 more `ldd`!!

Programs like this  
`your code` `z lib` `sqlite`  
are called "statically linked"

And programs like this  
`your code` `z lib` `sqlite`  
are called "dynamically linked"

Where the dynamic linker looks  
① DT\_RPATH in your executable  
② LD\_LIBRARY\_PATH  
③ DT\_RUNPATH in executable  
④ /etc/ld.so.cache  
(run ldconfig -p to see contents)  
⑤ /lib, /usr/lib

I got a "library not found" error when running my binary?!

If you know where the library is, try setting the LD\_LIBRARY\_PATH environment variable  
LD-LIBRARY-PATH tells me where to look!  
dynamic linker

# signals

7

If you've ever used `kill` you've used signals

DIE!!! `process`  
Okay! `process`

the Linux Kernel sends processes signals in lots of situations

- that pipe is closed
- your child terminated
- illegal instruction
- the timer you set expired
- segmentation fault

Every signal has a default action, which is one of:  
Ignore `process`  
Kill process `process`  
Kill process AND make core dump file `process`  
Stop process `process`  
Resume process `process`

You can send signals yourself with the kill system call or command

SIGINT	ctrl-C	various levels of
SIGTERM	Kill -q	"die"
SIGKILL		
SIGHUP	Kill -HUP	often interpreted as "reload config", e.g. by nginx

Signals can be hard to handle correctly since they can happen at ANY time

"oo handling a signal process  
SURPRISE! Another signal!"

# file descriptors

8

Unix systems use integers to track open files

`[!] Open foo.txt`

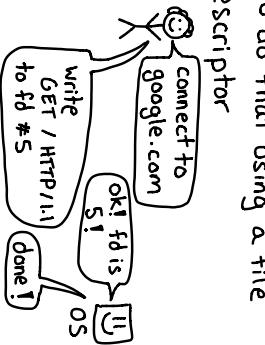
`[!] process`

`[!] kernel`

`[!] okay! That's file #7 for you.`

these integers are called file descriptors

When you read or write to a file / pipe / network connection you do that using a file descriptor



lsof (list open files) will show you a process's open files

`$ lsof -p 4242 ← PID we're interested in`

`FD NAME`

`1 /dev/pts/4tuy1`

`2 pipe: 29174`

`3 /home/bark/awesome.txt`

`5 /tmp/ ↗`

FD is for file descriptor

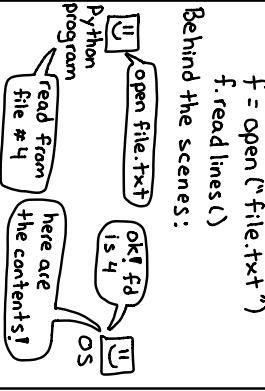
Let's see how some simple Python code works under the hood:

Python:

`f = open("file.txt")`

`f.readlines()`

Behind the scenes:



file descriptors can refer to:

→ files on disk

→ pipes

→ sockets (network connections)

→ terminals (like xterm)

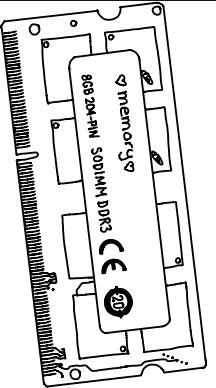
→ LOTS MORE (eventfd, inotify, signalfd, epoll, etc.)

`[!] not EVERYTHING on Unix is a file, but lots of things are`

# Virtual memory

17

your computer has physical memory



Linux keeps a mapping from virtual memory pages to physical memory pages called the page table

`[!] a "page" is a 4kb or bigger chunk of memory`

PID virtual addr physical addr

1971 0x20000 0x192000

2310 0x20000 0x228000

2310 0x21000 0x9788000

every program has its own virtual address space

`[!] 00 0x129520 → "puppies"`

`[!] program 1`

`[!] 00 0x129520 → "bananas"`

`[!] program 2`

every time you switch which process is running, Linux needs to switch the page table

`[!] here's the address of process 2950's page table`

`[!] Linux`

`[!] thanks, I'll use that now!`

`[!] MMU`

`[!] I'll look that up in the page table and then access the right physical address`

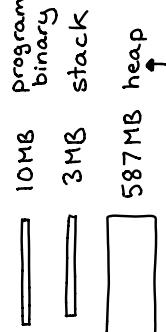
`[!] MMU management unit`

hardware

# memory allocation

16

Your program has memory



Your memory allocator (malloc) is responsible for 2 things.

THING 1: Keep track of what memory is used/free.



Your memory allocator's interface  
allocate size bytes of memory & return a pointer to it.  
malloc (size\_t size)  
free (void\* pointer)  
mark the memory as unused (and maybe give back to the OS).  
realloc(void\* pointer, size\_t size)  
ask for more/less memory for pointer:  
calloc (size\_t members, size\_t size)  
allocate array + initialize to 0.

THING 2 : Ask the OS for more memory!



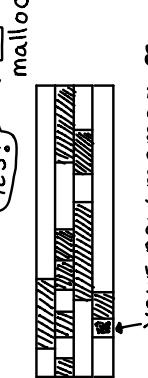
malloc isn't magic!!  
it's just a function!

You can always:

- use a different malloc library like jemalloc or tcmalloc (easy!)
- implement your own malloc (harder)

malloc tries to fill in unused space when you ask for memory

can I have 512 bytes of memory?  
your code



your new memory.

## pipes

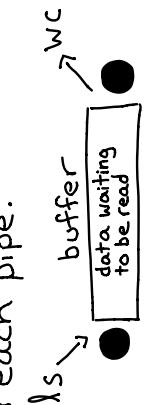
9

Sometimes you want to send the output of one process to the input of another

\$ ls | wc -l

53 ↴ 53 files!

Linux creates a buffer for each pipe.



If data gets written to the pipe faster than it's read, the buffer will fill up. When the buffer is full, writes to it will block (wait) until the reader reads. This is normal & ok.

when ls does  
write(●, "hi"),  
wc can read it!  
read(●)  
→ "hi"

Pipes are one way. → ●  
You can't write to ●.

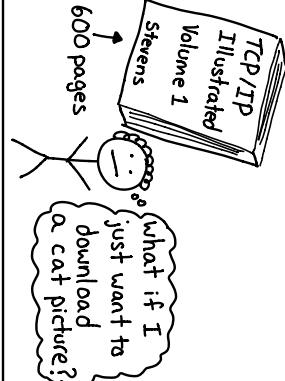
named pipes

\$ mkfifo my-pipe  
This lets 2 unrelated processes communicate through a pipe!  
f=open("./my-pipe")  
f.write("hi!\n")  
f=open("./my-pipe")  
f.readline() ← "hi!"

# Sockets

10

networking protocols are complicated



Every HTTP library uses sockets under the hood

```
$ curl awesome.com
```

Python: `requests.get("yay.us")`

oh, cool, I could write an HTTP library too if I wanted.\* Neat!

\* SO MANY edge cases though! ☹

Unix systems have an API called the "socket API" that makes it easier to make network connections

- ① Create a socket
- ② Connect to an IP/port ...
- ③ Write (fd, "GET /cat.png HTTP/1.1")
- ④ Read the response `cat-picture = read(fd ...)`

AF-INET?  
What's that?

AF-INET means basically "internet socket": it lets you connect to other computers on the internet using their IP address.

The main alternative is AF-UNIX ("unix domain socket") for connecting to programs on the same computer.

3 kinds of internet (AF-INET) sockets:

- SOCK\_STREAM = TCP  
curl uses this
- SOCK\_DGRAM = UDP  
dig (dns) uses this
- SOCK\_RAW = just let me send IP packets.  
ping uses this  
I will implement my own protocol.

## file buffering

15

I/O libraries don't always call write when you print.

`printf("I ❤ cats");`

`printf("I ❤ cats")`  
before actually writing

This is called buffering and it helps save on syscalls.

- when writing an interactive prompt!
- Python example:  
`print("password:", flush=True)`
- when you're writing to a pipe/socket
- no seriously, actually write to that pipe please

- ① None. This is the default for stderr.
- ② Line buffering.  
(write after newline). The default for terminals.
- ③ "full" buffering.  
(write in big chunks)  
The default for files and pipes.

## flushing



To force your I/O library to write everything it has in its buffer right now, call `flush()`!

`stdio`

I'll call write right away!!

- ① None. This is the default for stderr.
- ② Line buffering.  
(write after newline). The default for terminals.
- ③ "full" buffering.  
(write in big chunks)  
The default for files and pipes.

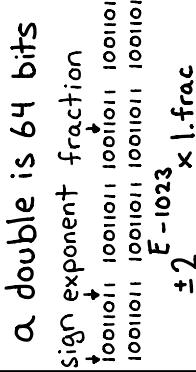
? ? ?  
I printed some text but it didn't appear on the screen. why??

time to learn about flushing!

?

# floating point

14

a double is 64 bits  
sign exponent fraction  


$E - 1023 \times 1.\text{frac}$   
 $\pm 2^{52}$   
That means there are  $2^{64}$  doubles.  
The biggest one is about  $2^{1023}$

weird double arithmetic

$$2^{52} + 0.2 = 2^{52} \quad \leftarrow \begin{array}{l} \text{(the next number after} \\ 2^{52} \text{ is } 2^{52+1}) \end{array}$$
$$1 + \frac{1}{2^{54}} = 1 \quad \leftarrow \begin{array}{l} \text{(the next number after} \\ 1 \text{ is } 1 + \frac{1}{2^{52}}) \end{array}$$
$$2^{2000} = \text{infinity} \quad \leftarrow \begin{array}{l} \text{infinity is a double} \\ \text{infinity - infinity} = \text{nan} \end{array}$$
$$\text{infinity - infinity} = \text{nan} \quad \leftarrow \begin{array}{l} \text{nan = "not a number"} \end{array}$$

doubles get farther apart as they get bigger between  $2^n$  and  $2^{n+1}$  there are always  $2^{52}$  doubles, evenly spaced.  
that means the next double after  $2^{60}$  is  $2^{60} + 64 = \frac{2^{60}}{2^{52}}$

JavaScript only has doubles (no integers!)

$$> 2^{**53}$$
  
$$900719254740992$$
$$> 2^{**53} + 1$$
  
$$900719254740992$$

same number! uh oh!

doubles are scary and their arithmetic is weird!

they're very logical!  
just understand how they work and don't use integers over  $2^{53}$  in JavaScript ♥

# unix domain sockets !!

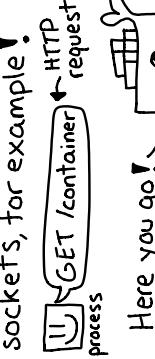
unix domain sockets are files.

\$ file mysock.sock  
socket

the file's permissions determine who can send data to the socket.

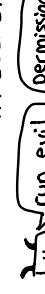
they let 2 programs on the same computer communicate.

Docker uses Unix domain sockets, for example!



Here you go!

advantage 1

Lets you use file permissions to restrict access to HTTP/ database services!  
chmod 600 secret.sock  
This is why Docker uses a unix domain socket.   
 run evil container  Linux

advantage 2

UDP sockets aren't always reliable (even on the same computer).  
unix domain datagram sockets are reliable!  
And they won't reorder packets!  
I can send data and I know it'll arrive



advantage 3

You can send a file descriptor over a unix domain socket.  
Useful when handling untrusted input files!

# What's in a process?

12

<b>PID</b>	<b>USER and GROUP</b>	<b>ENVIRONMENT VARIABLES</b>	<b>SIGNAL HANDLERS</b>
process #129 reporting for duty!	Who are you running as? julia!	like PATH! you can set them with: <code>\$ env A=val ./program</code>	I ignore SIGTERM! I shut down safely!
<b>WORKING DIRECTORY</b>	<b>PARENT PID</b>	<b>COMMAND LINE ARGUMENTS</b>	<b>OPEN FILES</b>
Relative paths ( <code>./blah</code> ) are relative to the working directory! chdir changes it.	PID 1 ( <code>init</code> ) is everyone's ancestor → PID 147 → PID 129	See them in <code>/proc/pid/cmdline</code>	Every open file has an offset.
<b>MEMORY</b>	<b>THREADS</b>	<b>CAPABILITIES</b>	<b>NAMESPACES</b>
heap! stack! Shared libraries! the program's binary! mmaped files!	sometimes one sometimes LOTS	I have <code>CAP_PTRACE</code> Well I have <code>CAP_SYS_ADMIN</code>	I'm in the host network namespace I have my own namespace! Container process

## threads

13

<b>Threads</b>	<b>Sharing memory</b>	<b>Why use threads</b>
Threads let a process do many different things at the same time	threads in the same process share memory	instead of starting a new process?
process:	and they share code <code>calculate-pi</code> <code>find-big-prime-number</code>	→ a thread takes less time to create. sharing data between threads is very easy. But it's also easier to make mistakes with threads.

**Thread 1**

I'm going to add 1 to that number!

**Thread 2**

I'm going to add 1 to that number!

**Memory**

23

same time

RESULT: 24 ← WRONG.  
Should be 25!

You weren't supposed to CHANGE that data!